

PATENT APPLICATION

POLYMORPHISM DETECTION

Inventors:

Robert J. Lipshutz, a citizen of
the United States, residing at:
970 Palo Alto Avenue
Palo Alto, CA 94301

Ronald Sapolisky, a citizen of
the United States, residing at:
630 Los Robles Avenue #16
Palo Alto, CA 94306

Ghassan Ghandour, a citizen of
the United States, residing at:
73 Palmer Lane
Atherton, CA 94027

Assignee:

Affymetrix, Inc.
3380 Central Expressway
Santa Clara, CA 95051

TOWNSEND and TOWNSEND and CREW LLP
Two Embarcadero Center, 8th Floor
San Francisco, California 94111-3834
(415) 576-0200

09939119-00822401

POLYMORPHISM DETECTION

5

July 1
This application is a continuation-in-part of 08/563,762, filed November 29, 1995, and claims the benefit of U.S. provisional application 60/017,260, filed May 10, 1996, the disclosures of which are incorporated by reference in their entirety for all purposes.

10

BACKGROUND OF THE INVENTION

00000

00005

00006

00007

00008

00009

00010

00011

00012

00013

00014

00015

00016

00017

00018

00019

00020

00021

00022

00023

00024

00025

00026

00027

00028

00029

00030

00031

00032

00033

00034

00035

The relationship between structure and function of macromolecules is of fundamental importance in the understanding of biological systems. These relationships are important to understanding, for example, the functions of enzymes, structural proteins and signaling proteins, ways in which cells communicate with each other, as well as mechanisms of cellular control and metabolic feedback.

Genetic information is critical in continuation of life processes. Life is substantially informationally based and its genetic content controls the growth and reproduction of the organism and its complements. The amino acid sequences of polypeptides, which are critical features of all living systems, are encoded by the genetic material of the cell. Further, the properties of these polypeptides, e.g., as enzymes, functional proteins, and structural proteins, are determined by the sequence of amino acids which make them up. As structure and function are integrally related, many biological functions may be explained by elucidating the underlying structural features which provide those functions, and these structures are determined by the underlying genetic information in the form of polynucleotide sequences. Further, in addition to encoding polypeptides, polynucleotide sequences also can be involved in control and regulation of gene expression. It therefore follows that the determination of the make-up of this genetic information has achieved significant scientific importance.

As a specific example, diagnosis and treatment of a variety of disorders may often be accomplished through identification and/or manipulation of the genetic material which encodes for specific disease associated traits. In order to accomplish this, however, one must first identify a correlation between a particular gene and a particular trait. This is generally accomplished by providing a genetic linkage map through which one identifies a set of genetic markers that follow a particular trait. These markers can identify the location of the gene encoding for that trait within the genome, eventually leading to the identification of the gene. Once the gene is identified, methods of treating the disorder that result from that gene, i.e., as a result of overexpression, constitutive expression, mutation, underexpression, etc., can be more easily developed.

One class of genetic markers includes variants in the genetic code termed "polymorphisms." In the course of evolution, the genome of a species can collect a number of variations in individual bases. These single base changes are termed single-base polymorphisms. Polymorphisms may also exist as stretches of repeating sequences that vary as to the length of the repeat from individual to individual. Where these variations are recurring, e.g., exist in a significant percentage of a population, they can be readily used as markers linked to genes involved in mono- and polygenic traits. In the human genome, single-base polymorphisms occur roughly once per 300 bp. Though many of these variant bases appear too infrequently among the allele population for use as genetic markers (i.e., <1%), useful polymorphisms (e.g., those occurring in 20 to 50 % of the allele population) can be found approximately once per kilobase. Accordingly, in a human genome of approximately 3 Gb, one would expect to find approximately 3,000,000 of these "useful" polymorphisms.

The use of polymorphisms as genetic linkage markers is thus of critical importance in locating, identifying and characterizing the genes which are responsible for specific traits. In particular, such mapping techniques allow for the identification of genes responsible for a variety of disease

or disorder-related traits which may be used in the diagnosis and or eventual treatment of those disorders. Given the size of the human genome, as well as those of other mammals, it would generally be desirable to provide methods of rapidly 5 identifying and screening for polymorphic genetic markers. The present invention meets these and other needs.

SUMMARY OF THE INVENTION

One aspect of the invention is an array of 10 oligonucleotide probes for detecting a polymorphism in a target nucleic acid sequence using Principal Component Analysis, said array comprising at least one detection block of probes, said detection block including a first group of probes that are complementary to said target nucleic acid sequence except that the group of probes includes all possible monosubstitutions of positions in said sequence that are within n bases of a base in said sequence that is 15 complementary to said polymorphism, wherein n is from 0 to 5, and a second and third group of probes complementary to marker-specific regions upstream and downstream of the target nucleic acid sequence, wherein the third group of probes differs from the second set of probes at single bases 20 corresponding to known mismatch positions.

A further aspect of the invention is a method of 25 identifying whether a target nucleic acid sequence includes a polymorphic variant using principal component analysis, comprising:

hybridizing said target nucleic acid sequence to said array comprising at least one detection block of 30 probes, said detection block including a first group of probes that are complementary to said target nucleic acid sequence except that the group of probes includes all possible monosubstitutions of positions in said sequence that are within n bases of a base in said sequence that is 35 complementary to said polymorphism, wherein n is from 0 to 5, and a second and third group of probes complementary to marker-specific regions upstream and downstream of the target nucleic acid sequence, wherein the third group of probes

differs from the second set of probes at single bases corresponding to known mismatch positions; and

determining hybridization intensities of the target nucleic acid and the marker-specific regions to identify said polymorphic variant. In one embodiment of the invention, the step of determining comprises:

a) calculating the control difference between the average of the hybridization intensities of the second group of probes, the hybridization intensities comprising control perfect matches (PM), minus the average of the hybridization intensities, the hybridization intensities comprising control single-base mismatches (MM);

b) calculating the possible perfect match intensity and a heteromismatch intensity from the hybridization intensities for each position of monosubstitutions of the first group of probes;

c) calculating the difference between the possible perfect match intensity and the heteromismatch intensity for each position of monosubstitutions of the first group of probes;

d) calculating a normalized difference (ND) by dividing the difference of step (c) by the control difference;

(e) using principal component analysis, identifying a polymorphism by comparing normalized differences between individuals in a population.

BRIEF DESCRIPTION OF THE DRAWINGS

Figure 1 shows a schematic illustration of light-directed synthesis of oligonucleotide arrays.

Figure 2A shows a schematic representation of a single oligonucleotide array containing 78 separate detection blocks. Figure 2B shows a schematic illustration of a detection block for a specific polymorphism denoted WI-567. Figure 2B also shows the triplet layout of detection blocks for the polymorphism employing 20-mer oligonucleotide probes having substitutions 7, 10 and 13 bp from the 3' end of the probe. The probes present in the shaded portions of each detection block are shown adjacent to each detection block.

5
567A4
Figure 3 illustrates a tiling strategy for a polymorphism denoted WI-567, and having the sequence 5'-TGCTGCCTTGGTTC[A/G]AGCCCTCATCTCTTT-3'. A detection block specific for the WI-567 polymorphism is shown with the probe sequences tiled therein listed above. Predicted patterns for both homozygous forms and the heterozygous form are shown at the bottom.

10
567A5
Figure 4 shows a schematic representation of a detection block specific for the polymorphism denoted WI-1959 having the sequence 5'-ACCAAAATCAGTC[T/C]GGGTAAGTGAGAGTG-3' with the polymorphism indicated by the brackets. A fluorescent scan of hybridization of the heterozygous and both homozygous forms are shown in the center, with the predicted hybridization pattern for each being indicated below.

15
00116660
Figure 5 illustrates an example of a computer system used to execute the software of the present invention which determines whether polymorphic markers in DNA are heterozygote, homozygote with a first polymorphic marker or homozygote with a second polymorphic marker.

20
002200
Figure 6 shows a system block diagram of computer system 1 used to execute the software of the present invention.

25
0042200
Figure 7 shows a probe array including probes with base substitutions at base positions within two base positions of the polymorphic marker. The position of the polymorphic marker is denoted P_0 and which may have one of two polymorphic markers x and y (where x and y are one of A, C, G, or T).

30
0042200
Figure 8 shows a probe array including probes with base substitutions at base positions within two base positions of the polymorphic marker.

Figure 9 shows a high level flowchart of analyzing intensities to determine whether polymorphic markers in DNA are heterozygote, homozygote with a first polymorphic marker or homozygote with a second polymorphic marker.

35
0042200
Figure 10 shows a Principal Components Plot of Marker 219 (KRT8m1).

Figure 11 shows a schematic representation of a process for carrying out the polymorphism detection methods of the invention.

Figure 12 shows the algorithms used for identifying 5 genotypes, using the methods of the present invention.

Figure 13 shows the DB scores of one marker plotted along with the genotypes determined by standard sequencing. Approximately 220 biallelic markers were assayed together for each individual for a sixteen member family.

10

DETAILED DESCRIPTION OF THE INVENTION

I. General

The present invention generally provides rapid and efficient methods for screening samples of genomic material for polymorphisms, and arrays specifically designed for carrying out these analyses. In particular, the present invention relates to the identification and screening of single base polymorphisms in a sample. In general, the methods of the present invention employ arrays of oligonucleotide probes that are complementary to target nucleic acid sequence segments from an individual (e.g., a human or other mammal) which target sequences include specific identified polymorphisms, or "polymorphic markers." The probes are typically arranged in detection blocks, each block being capable of discriminating the three genotypes for a given marker, e.g., the heterozygote or either of the two homozygotes. The method allows for rapid, automatable analysis of genetic linkage to even complex polygenic traits.

Oligonucleotide arrays typically comprise a plurality of different oligonucleotide probes that are coupled to a surface of a substrate in different known locations. These oligonucleotide arrays, also described as "Genechips™," have been generally described in the art, for example, U.S. Patent No. 5,143,854 and PCT patent publication Nos. WO 90/15070 and 92/10092. These arrays may generally be produced using mechanical synthesis methods or light directed synthesis methods which incorporate a combination of photolithographic methods and solid phase oligonucleotide synthesis methods.

See Fodor et al., *Science*, 251:767-777 (1991), Pirrung et al., U.S. Patent No. 5,143,854 (see also PCT Application No. WO 90/15070) and Fodor et al., PCT Publication No. WO 92/10092 and U.S. Patent No. 5,424,186, each of which is hereby incorporated herein by reference. Techniques for the synthesis of these arrays using mechanical synthesis methods are described in, e.g., U.S. Patent No. 5,384,261, incorporated herein by reference in its entirety for all purposes.

The basic strategy for light directed synthesis of oligonucleotides on a VLSIPS™ Array is outlined in Figure 1. The surface of a substrate or solid support, modified with photosensitive protecting groups (X) is illuminated through a photolithographic mask, yielding reactive hydroxyl groups in the illuminated regions. A selected nucleotide, typically in the form of a 3'-O-phosphoramidite-activated deoxynucleoside (protected at the 5' hydroxyl with a photosensitive protecting group), is then presented to the surface and coupling occurs at the sites that were exposed to light. Following capping and oxidation, the substrate is rinsed and the surface is illuminated through a second mask, to expose additional hydroxyl groups for coupling. A second selected nucleotide (e.g., 5'-protected, 3'-O-phosphoramidite-activated deoxynucleoside) is presented to the surface. The selective deprotection and coupling cycles are repeated until the desired set of products is obtained. Pease et al., *Proc. Natl. Acad. Sci.* (1994) 91:5022-5026. Since photolithography is used, the process can be readily miniaturized to generate high density arrays of oligonucleotide probes. Furthermore, the sequence of the oligonucleotides at each site is known.

II. Identification of Polymorphisms

The methods and arrays of the present invention primarily find use in the identification of so-called "useful" (i.e., those that are present in approximately 20% or more of the allele population). The present invention also relates to the detection or screening of specific variants of previously identified polymorphisms.

A wide variety of methods can be used to identify specific polymorphisms. For example, repeated sequencing of genomic material from large numbers of individuals, although extremely time consuming, can be used to identify such polymorphisms. Alternatively, ligation methods may be used, where a probe having an overhang of defined sequence is ligated to a target nucleotide sequence derived from a number of individuals. Differences in the ability of the probe to ligate to the target can reflect polymorphisms within the sequence. Similarly, restriction patterns generated from treating a target nucleic acid with a prescribed restriction enzyme or set of restriction enzymes can be used to identify polymorphisms. Specifically, a polymorphism may result in the presence of a restriction site in one variant but not in another. This yields a difference in restriction patterns for the two variants, and thereby identifies a polymorphism. In a related method, U.S. Patent Application Serial No. 08/485,606, filed June 7, 1995 describes a method of identifying polymorphisms using type-IIIs endonucleases to capture ambiguous base sequences adjacent the restriction sites, and characterizing the captured sequences on oligonucleotide arrays. The patterns of these captured sequences are compared from various individuals, the differences being indicative of potential polymorphisms.

In a preferred aspect, the identification of polymorphisms takes into account the assumption that a useful polymorphism (i.e., one that occurs in 20 to 50% of the allele population) occurs approximately once per 1 KB in a given genome. In particular, random sequences of a genome, e.g., random 1 kb sequences of the human genome such as expressed sequence tags or "ESTs", can be sequenced from a limited number of individuals. When a variant base is detected with sufficient frequency, it is designated a "useful" polymorphism. In practice, the method generally analyzes the same 1 KB sequence from a small number of unrelated individuals, i.e., from 3 to 5 (6 to 10 alleles). Where a variant sequence is identified, it is then compared to a separate pool of material from unrelated individuals (i.e., 10

unrelated individuals). Where the variant sequence identified from the first set of individuals is detectable in the pool of the second set, it is assumed to exist at a sufficiently high frequency, e.g., at least about 20% of the allele population, 5 thereby qualifying as a useful marker for genetic linkage analysis.

III. Screening Polymorphisms

Screening polymorphisms in samples of genomic material according to the methods of the present invention, is generally carried out using arrays of oligonucleotide probes. These arrays may generally be "tiled" for a large number of specific polymorphisms. By "tiling" is generally meant the synthesis of a defined set of oligonucleotide probes which is made up of a sequence complementary to the target sequence of interest, as well as preselected variations of that sequence, e.g., substitution of one or more given positions with one or more members of the basis set of monomers, i.e. nucleotides. Tiling strategies are discussed in detail in Published PCT Application No. WO 95/11995, incorporated herein by reference in its entirety for all purposes. By "target sequence" is meant a sequence which has been identified as containing a polymorphism, and more particularly, a single-base polymorphism, also referred to as a "biallelic base." It will be understood that the term "target sequence" is intended to encompass the various forms present in a particular sample of genomic material, i.e., both alleles in a diploid genome.

In a particular aspect, arrays are tiled for a number of specific, identified polymorphic marker sequences. In particular, the array is tiled to include a number of detection blocks, each detection block being specific for a specific polymorphic marker or set of polymorphic markers. For example, a detection block may be tiled to include a number of probes which span the sequence segment that includes a specific polymorphism. To ensure probes that are complementary to each variant, the probes are synthesized in pairs differing at the biallelic base.

In addition to the probes differing at the biallelic bases, monosubstituted probes are also generally tiled within the detection block. These monosubstituted probes have bases at and up to a certain number of bases in either direction from the polymorphism, substituted with the remaining nucleotides (selected from A, T, G, C or U). Typically, the probes in a tiled detection block will include substitutions of the sequence positions up to and including those that are 5 bases away from the base that corresponds to the polymorphism. Preferably, bases up to and including those in positions 2 bases from the polymorphism will be substituted. The monosubstituted probes provide internal controls for the tiled array, to distinguish actual hybridization from artifactual cross-hybridization. An example of this preferred substitution pattern is shown in Figure 3.

A variety of tiling configurations may also be employed to ensure optimal discrimination of perfectly hybridizing probes. For example, a detection block may be tiled to provide probes having optimal hybridization intensities with minimal cross-hybridization. For example, where a sequence downstream from a polymorphic base is G-C rich, it could potentially give rise to a higher level of cross-hybridization or "noise," when analyzed. Accordingly, one can tile the detection block to take advantage of more of the upstream sequence. Such alternate tiling configurations are schematically illustrated in Figure 2B, bottom, where the base in the probe that is complementary to the polymorphism is placed at different positions in the sequence of the probe relative to the 3' end of the probe. For ease of discussion, both the base which represents the polymorphism and the complementary base in the probe are referred to herein as the "polymorphic base" or "polymorphic marker."

Optimal tiling configurations may be determined for any particular polymorphism by comparative analysis. For example, triplet or larger detection blocks like those illustrated in Figure 2B may be readily employed to select such optimal tiling strategies.

Additionally, arrays will generally be tiled to provide for ease of reading and analysis. For example, the probes tiled within a detection block will generally be arranged so that reading across a detection block the probes are tiled in succession, i.e., progressing along the target sequence one or more base at a time (See, e.g., Figure 3, middle).

Once an array is appropriately tiled for a given polymorphism or set of polymorphisms, the target nucleic acid is hybridized with the array and scanned. Hybridization and scanning are generally carried out by methods described in, e.g., Published PCT Application Nos. WO 92/10092 and WO 95/11995, and U.S. Patent No. 5,424,186, previously incorporated herein by reference in their entirety for all purposes. In brief, a target nucleic acid sequence which includes one or more previously identified polymorphic markers is amplified by well known amplification techniques, e.g., PCR. Typically, this involves the use of primer sequences that are complementary to the two strands of the target sequence both upstream and downstream from the polymorphism. Asymmetric PCR techniques may also be used. Amplified target, generally incorporating a label, is then hybridized with the array under appropriate conditions. Upon completion of hybridization and washing of the array, the array is scanned to determine the position on the array to which the target sequence hybridizes. The hybridization data obtained from the scan is typically in the form of fluorescence intensities as a function of location on the array.

Although primarily described in terms of a single detection block, e.g., for detection of a single polymorphism, in preferred aspects, the arrays of the invention will include multiple detection blocks, and thus be capable of analyzing multiple, specific polymorphisms. For example, preferred arrays will generally include from about 50 to about 4000 different detection blocks with particularly preferred arrays including from 100 to 3000 different detection blocks.

In alternate arrangements, it will generally be understood that detection blocks may be grouped within a

single array or in multiple, separate arrays so that varying, optimal conditions may be used during the hybridization of the target to the array. For example, it may often be desirable to provide for the detection of those polymorphisms that fall within G-C rich stretches of a genomic sequence, separately from those falling in A-T rich segments. This allows for the separate optimization of hybridization conditions for each situation.

10 IV. Calling

After hybridization and scanning, the hybridization data from the scanned array is then analyzed to identify which variant or variants of the polymorphic marker are present in the sample, or target sequence, as determined from the probes to which the target hybridized, e.g., one of the two homozygote forms or the heterozygote form. This determination is termed "calling" the genotype. Calling the genotype is typically a matter of comparing the hybridization data for each potential variant, and based upon that comparison, identifying the actual variant (for homozygotes) or variants (for heterozygotes) that are present. In one aspect, this comparison involves taking the ratio of hybridization intensities (corrected for average background levels) for the expected perfectly hybridizing probes for a first variant versus that of the second variant. Where the marker is homozygous for the first variant, this ratio will be a large number, theoretically approaching an infinite value. Where homozygous for the second variant, the ratio will be a very low number, i.e., theoretically approaching zero. Where the marker is heterozygous, the ratio will be approximately 1. These numbers are, as described, theoretical. Typically, the first ratio will be well in excess of 1, i.e., 2, 4, 5 or greater. Similarly, the second ratio will typically be substantially less than 1, i.e., 0.5, 0.2, 0.1 or less. The ratio for heterozygotes will typically be approximately equal to 1, i.e. from 0.7 to 1.5. These ratios can vary based upon the specific sequence surrounding the polymorphism, and can

also be adjusted based upon a standard hybridization with a control sample containing the variants of the polymorphism.

The quality of a given call for a particular genotype may also be checked. For example, the maximum perfect match intensity can be divided by a measure of the background noise (which may be represented by the standard deviation of the mismatched intensities). Where the ratio exceeds some preselected cut-off point, the call is determined to be good. For example, where the maximum intensity of the expected perfect matches exceeds twice the noise level, it might be termed a good call. In an additional aspect, the present invention provides software for performing the above described comparisons.

Fig. 5 illustrates an example of a computer system used to execute the software of the present invention which determines whether polymorphic markers in DNA are heterozygote, homozygote with a first variant of a polymorphism or homozygote with a second variant of a polymorphism. Fig. 5 shows a computer system 1 which includes a monitor 3, screen 5, cabinet 7, keyboard 9, and mouse 11. Mouse 11 may have one or more buttons such as mouse buttons 13. Cabinet 7 houses a CD-ROM drive 15 or a hard drive (not shown) which may be utilized to store and retrieve software programs incorporating the present invention, digital images for use with the present invention, and the like. Although a CD-ROM 17 is shown as the removable media, other removable tangible media including floppy disks, tape, and flash memory may be utilized. Cabinet 7 also houses familiar computer components (not shown) such as a processor, memory, and the like.

Fig. 6 shows a system block diagram of computer system 1 used to execute the software of the present invention. As in Fig. 5, computer system 1 includes monitor 3 and keyboard 9. Computer system 1 further includes subsystems such as a central processor 102, system memory 104, I/O controller 106, display adapter 108, removable disk 112, fixed disk 116, network interface 118, and speaker 120. Other computer systems suitable for use with the present invention

may include additional or fewer subsystems. For example, another computer system could include more than one processor 102 (i.e., a multi-processor system) or a cache memory.

Arrows such as 122 represent the system bus

5 architecture of computer system 1. However, these arrows are illustrative of any interconnection scheme serving to link the subsystems. For example, a local bus could be utilized to connect the central processor to the system memory and display adapter. Computer system 1 shown in Fig. 6 is but an example 10 of a computer system suitable for use with the present invention. Other configurations of subsystems suitable for use with the present invention will be readily apparent to one of ordinary skill in the art.

Fig. 7 shows a probe array including probes with base substitutions at base positions within two base positions of the polymorphic marker. The position of the polymorphic marker is denoted P_0 and which may have one of two variants of the polymorphic markers x and y (where x and y are one of A, C, G, or T). As indicated, at P_{-2} there are two columns of four cells which contain a base substitution two base positions to the left, or 3', from the polymorphic marker. The column denoted by an "x" contains polymorphic marker x and the column denoted by a "y" contains polymorphic marker y.

Similarly, P_{-1} contains probes with base substitutions one base position to the left, or 3', of the polymorphic marker. P_0 contains probes with base substitutions at the polymorphic marker position. Accordingly, the two columns in P_0 are identical. P_1 and P_2 contain base substitutions one and two base positions to the 30 right, or 5', of the polymorphic marker, respectively.

As a hypothetical example, assume a single base polymorphism exists where one allele contains the subsequence TCAAG whereas another allele contains the subsequence TCGAG, where the underlined base indicates the polymorphism in each 35 allele. Fig. 8 shows a probe array including probes with base substitutions at base positions within two base positions of the polymorphic marker. In the first two columns, the cells which contain probes with base A (complementary to T in the

0005660
5554445
4443335
3332225
2221115
1110005
0009990

alleles) two positions from the left of the polymorphic marker are shaded. They are shaded to indicate that it is expected that these cells would exhibit the highest hybridization to the labeled sample nucleic acid. Similarly, the second two 5 columns have cells shaded which have probes with base G (complementary to C in the alleles) one position to the left of the polymorphic marker.

At the polymorphic marker position (corresponding to P_0 in Fig. 7), there are two columns: one denoted by an "A" 10 and one denoted by a "G". Although, as indicated earlier, the probes in these two columns are identical, the probes contain base substitutions for the polymorphic marker position. An "N" indicates the cells that have probes which are expected to exhibit a strong hybridization if the allele contains a 15 polymorphic marker A. As will become apparent in the following paragraphs, "N" stands for numerator because the intensity of these cells will be utilized in the numerator of an equation. Thus, the labels were chosen to aid the reader's understanding of the present invention.

A "D" indicates the cells that have probes which are expected to exhibit a strong hybridization if the allele contains a polymorphic marker G. "D" stands for denominator because the intensity of these cells will be utilized in the 20 denominator of an equation. The "n" and "d" labeled cells indicate these cells contain probes with a single base mismatch near the polymorphic marker. As before, the labels indicate where the intensity of these cells will be utilized in a following equation.

Fig. 9 shows a high level flowchart of analyzing 25 intensities to determine whether polymorphic markers in DNA are heterozygote, homozygote with a first polymorphic marker or homozygote with a second polymorphic marker. At step 202, the system receives the fluorescent intensities of the cells on the chip. Although in a preferred embodiment, the 30 hybridization of the probes to the sample are determined from fluorescent intensities, other methods and labels including radioactive labels may be utilized with the present invention. An example of one embodiment of a software program for

carrying out this analysis is reprinted in Software Appendix A.

A perfect match (PM) average for a polymorphic marker x is determined by averaging the intensity of the cells at P_0 that have the base substitution equal to x in Fig. 7. Thus, for the example in Fig. 8, the perfect match average for A would add the intensities of the cells denoted by "N" and divide the sum by 2.

A mismatch (MM) average for a polymorphic marker x is determined by averaging the intensity of the cells that contain the polymorphic marker x and a single base mismatch in Fig. 7. Thus, for the example in Fig. 8, the mismatch average for A would be the sum of cells denoted by "n" and dividing the sum by 14.

A perfect match average and mismatch average for polymorphic marker y is determined in a similar manner utilizing the cells denoted by "D" and "d", respectively. Therefore, the perfect match averages are an average intensity of cells containing probes that are perfectly complementary to an allele. The mismatch averages are an average of intensity of cells containing probes that have a single base mismatch near the polymorphic marker in an allele.

At step 204, the system calculates a Ratio of the perfect match and mismatch averages for x to the perfect match and mismatch averages for y . The numerator of the Ratio includes the mismatch average for x subtracted from the perfect mismatch for x . In a preferred embodiment, if the resulting numerator is less than 0, the numerator is set equal to 0.

The denominator of the Ratio includes the mismatch average for y subtracted from the perfect mismatch for y . In a preferred embodiment, if the resulting denominator is less than or equal to 0, the numerator is set equal to a de minimum value like 0.00001.

Once the system has calculated the Ratio, the system calculates DB at step 206. DB is calculated by the equation $DB = 10 * \log_{10} \text{Ratio}$. The logarithmic function puts the ratio on

a linear scale and makes it easier to interpret the results of the comparison of intensities.

At step 208, the system performs a statistical check on the data or hybridization intensities. The statistical check is performed to determine if the data will likely produce good results. In a preferred embodiment, the statistical check involves testing whether the maximum of the perfect match averages for x or y is at least twice as great as the standard deviation of the intensities of all the cells containing a single base mismatch (i.e., denoted by a "n" or "d" in Fig. 8). If the perfect match average is at least two times greater than this standard deviation, the data is likely to produce good results and this is communicated to the user.

The system analyzes DB at step 210 to determine if DB is approaching $-\infty$, near 0, or approaching $+\infty$. If DB is approaching a negative infinity, the system determines that the sample DNA contains a homozygote with a first polymorphic marker corresponding to x at step 212. If DB is near 0, the system determines that the sample DNA contains a heterozygote corresponding to both polymorphic markers x and y at step 214. Although described as approaching ∞ , etc., as described previously, these numbers will generally vary, but are nonetheless indicative of the calls described. If DB is approaching a positive infinity, the system determines that the sample DNA contains a homozygote with a second polymorphic marker corresponding to y at step 216.

A visual inspection of the Ratio equation in step 204 shows that the numerator should be higher than the denominator if the DNA sample only has the polymorphic marker corresponding to x. Similarly, the denominator should be higher than the numerator if the DNA sample only has a polymorphic marker corresponding to y. If the DNA sample has both polymorphic markers, indicating a heterozygote, the Ratio should be approximately equal to 1 which results in a 0 when the logarithm of the Ratio is calculated.

The equations discussed above illustrate just one embodiment of the present invention. These equations have correctly identified polymorphic markers when a visual

inspection would seem to indicate a different result. This may be the case because the equations take into account the mismatch intensities in order to determine the presence or absence of the polymorphic markers.

Those of skill in the art, upon reading the instant disclosure will appreciate that a variety of modifications to the above described methods and arrays may be made without departing from the spirit or scope of the invention. For

example, one may select the strand of the target sequence to optimize the ability to call a particular genome.

Alternatively, one may analyze both strands, in parallel, to provide greater amounts of data from which a call can be made. Additionally, the analyses, i.e., amplification and scanning may be performed using DNA, RNA, mixed polymers, and the like.

The present invention is further illustrated by the following examples. These examples are merely to illustrate aspects of the present invention and are not intended as limitations of this invention.

v. Examples

Example 1- Chip Tiling

A DNA chip is prepared which contains three detection blocks for each of 78 identified single base polymorphisms or biallelic markers, in a segment of human DNA (the "target" nucleic acid). Each detection block contains probes wherein the identified polymorphism occurs at the position in the target nucleic acid complementary to the 7th, 10th and 13th positions from the 3' end of 20-mer oligonucleotide probes. A schematic representation of a single oligonucleotide array containing all 78 detection blocks is shown in Figure 2A.

The tiling strategy for each block substitutes bases in the positions at, and up to two bases, in either direction from the polymorphism. In addition to the substituted positions, the oligonucleotides are synthesized in pairs differing at the biallelic base. Thus, the layout of the detection block (containing 40 different oligonucleotide probes) allows for controlled comparison of the sequences

involved, as well as simple readout without need for complicated instrumentation. A schematic illustration of this tiling strategy within a single detection block is shown in Figure 3, for a specific polymorphic marker denoted WI-567.

5 Example 2- Detection of Polymorphisms

A target nucleic acid is generated from PCR products amplified by primers flanking the markers. These amplicons can be produced singly or in multiplexed reactions. Target can be produced as ss-DNA by asymmetric PCR from one primer 10 flanking the polymorphism or as RNA transcribed in vitro from promoters linked to the primers. Fluorescent label is introduced into target directly as dye-bearing nucleotides, or bound after amplification using dye-streptavidin complexes to incorporated biotin containing nucleotides. In DNA produced by asymmetric PCR fluorescent dye is linked directly to the 5' end of the primer.

Hybridization of target to the arrays tiled in Example 1, and subsequent washing are carried out with standard solutions of salt (SSPE, TMACl) and nonionic detergent (Triton-X100), with or without added organic solvent (formamide). Targets and markers generating strong signals are washed under stringent hybridization conditions (37-40°C; 10% formamide; 0.25xSSPE washes) to give highly discriminating detection of the genotype. Markers giving lower hybridization 25 intensity are washed under less stringent conditions (\leq 30°C; 3M TMACl, or 6xSSPE; 6x and 1x SSPE washes) to yield highly discriminating detection of the genotype.

Detection of one polymorphic marker is illustrated in Figure 3. Specifically, a typical detection block is shown for the polymorphism denoted WI-1959, having the sequence 5'-ACCAAAAAATCAGTC[T/C]GGGTAACTGAGAGTG-3' with the polymorphism indicated by the brackets (Figure 3, top), for which all three genotypes are available (T/C heterozygote, C/C homozygote and T/T homozygote). The expected hybridization pattern for the homozygote and heterozygote targets are shown in Figure 3, bottom. Three chips were tiled with each chip including the illustrated detection block. Each block contained probes 35 having the substituted bases at the 7th, 10th and 13th

T04280-6TT6E660
00250
T01

positions from the 3' end of 20-mer oligonucleotide probes (20/7, 20/10 and 20/13, respectively). These alternate detection blocks were tiled to provide a variety of sequences flanking the polymorphism itself, to ensure at least one 5 detection block hybridizing with a sufficiently low background intensity for adequate detection.

Fluorouracil containing RNA was synthesized from a T7 promoter on the upstream primer, hybridized to the detection array in 6xSSPE + Triton-X100 at 30°C, and washed in 10 0.25xSSPE at room temperature. As shown in the scan Figure 3, middle, fluorescent scans of the arrays correctly identified the 5 homozygote or 10 heterozygote features.

Example 3- Alternate Gene Calling Method

An alternate method for calling the genotypes of a pedigree (or any collection of individuals) from P246 chip data is described herein. In particular, each sample from each of the individuals studied is amplified and hybridized to a P246 chip. The 246 chip employs a poly-tiling scheme and contains marker-specific control probes covering regions upstream and downstream from the single-base polymorphism. The significance of the control probes is that given that the target sample is amplified, these probes will display both perfect matches (PM) and single-base mismatches (MM) at known mismatch positions regardless of the target genotype. Even though in this document our description of the genotype calling method is based on the intensity data of a specific offset block (7/20, 10/20, 13/20) and a specific strand (T7, T3), this method is easily generalizable to accommodate data 25 combining multiple offset blocks and strands.

Considering the data for a collection of individuals for a given marker and a given offset block and strand, the relevant raw data for each individual are (1) two control PM intensities, (2) the corresponding (mismatch position = offset) control MM intensities, (3) 40 block intensities, (4) interrogations for each of 2 alleles for each of 5 positions. The averages of the two control PMs and of the two control MMs are computed. The difference of the two averages (PM - MM) is

labeled the Control Difference. For each of the 10 sets of 4 intensities (5 positions for each allele), the "possible" PM intensity is identified. (Note that in individuals where the allele is present a PM will occur in a predetermined probe from the set of 4 interrogation probes, that probe is what we call "possible" PM.) The hetero-MM probe (the same nucleotide appears at the mismatch position in both strands) is selected from the remaining 3, and the PM-MM difference is calculated. Each of the 10 block PM-MM values is divided by the Control Difference giving a Normalized Difference (ND). Thus the data for each individual are now reduced from 44 values to 10 ND values, 5 for each allele.

Two principal components analyses (PCA) are then performed (see Figure 10); one for each allele. PCA methodology originated with K. Pearson (1901) *Philosophical Magazine*, 2, 559-572 as a means of fitting planes to data by orthogonal least squares, but was later proposed by Hotelling (1933) *Journal of Educational Psychology*, 26:417-441, 498-520 and Hotelling (1936) *Psychometrika*, 1:27-35. for the particular purpose of analyzing correlation structure. The correlation structure analyzed in our case is that of the 5 ND values correlated over individuals. PCA attempts to find hierarchical sets of coefficients so that the simple weighted average of the ND values using the first set of coefficients would account for the largest portion of the variability among individuals. The second set would account for the largest portion of remaining variability with the constraint that it is orthogonal (non-overlapping) to the first, and so on. Without being bound to a particular theory, it is believed that the major source of variability among individuals on the 5 ND scores is mainly due to the difference between those with the given allele and those without. Thus, it can be expected that the first principal component would capture this difference, so that the weighted average based on the first set of coefficients is computed, individuals with the allele will generally have high scores and those without will have low scores. Moreover, combining the PCA results from the two alleles will distinguish between homozygous individuals with

the first allele (high PC scores on first low on second), homozygous individuals with the second allele (high PC scores on second low on first) and heterozygous individuals who will be high on both. This is illustrated in the enclosed plot of the two principal components for a set of 16 individuals from the K104 CEPH family. Clearly the individuals can be divided into three groups corresponding to the indicated genotypes.

10 Among the 221 biallelic markers detected on the polymorphism chip, seventeen of these markers were selected and assayed in a sixteen member CEPH family for genotyping by the present methods (see Figure 14) and by ABI sequencing. Of the 272 genotypes called between the two methods, there were only three disagreements (~90% concordance). All the genotypes called by either method were consistent with 15 Mendelian inheritance.

TDH280-611660

While the foregoing invention has been described in some detail for purposes of clarity and understanding, it will be clear to one skilled in the art from a reading of this disclosure that various changes in form and detail can be made without departing from the true scope of the invention. All publications and patent documents cited in this application are incorporated by reference in their entirety for all purposes to the same extent as if each individual publication or patent document were so individually denoted.

APPENDIX A

SOFTWARE APPENDIX

• **fullcal.awk** takes input from a POLYCHIP CCL file (1115 x 130) and extracts ratio information for every block on the chip. 22

101082419-003919

```

name(8,2) = "AGTCGC"
name(9,2) = "AARS"
name(9,3) = "TGATAT"
name(0,4) = "TTTA"
name(0,5) = "CTTCCC"
name(1,4) = "ATTC"
name(1,5) = "GCACAT"
name(2,4) = "BCL2"
name(2,5) = "ACGAGG"
name(3,4) = "BRCAGA"
name(3,5) = "CATGTC"
name(4,4) = "AGGAGC"
name(4,5) = "AGGAGA"
name(5,4) = "BNCAGC"
name(5,5) = "GAGAGC"
name(6,4) = "D382"
name(6,5) = "CCAGGT"
name(7,4) = "D3811"
name(7,5) = "TCTGAA"
name(8,4) = "D3812"
name(8,5) = "CCAGGT"
name(9,4) = "D3813"
name(9,5) = "CACTGC"
name(0,5) = "GCCACT"
name(1,5) = "GCK"
name(1,6) = "GAGACA"
name(2,5) = "HTT"
name(2,6) = "TCTTAC"
name(3,5) = "TTC"
name(3,6) = "TCTAAC"
name(4,5) = "HTT"
name(4,6) = "ACTTCA"
name(5,3) = "TGF2"
name(5,5) = "GCCACT"
name(6,5) = "TGTGAA"
name(6,6) = "TGTGAA"
name(7,5) = "TTC"
name(7,6) = "TCTTAC"
name(8,5) = "LDEA"
name(8,6) = "GCTCAA"
name(9,3) = "LTF2"
name(9,5) = "CCAGGG"
name(0,6) = "LTF"
name(0,7) = "CCAGGG"
name(1,6) = "HCC"
name(2,4) = "GCTGCA"
name(2,5) = "HTHT"
name(2,6) = "CCCTGG"
name(3,4) = "HAGACG"
name(3,6) = "CAGATG"
name(4,6) = "PAR"
name(4,7) = "PAR"
name(5,6) = "PAR/ADP"
name(5,7) = "GAGGAA"
name(6,5) = "GAGGAA"
name(6,6) = "PPPP3RL"
name(6,7) = "GACTAA"
name(7,6) = "RDR"
name(7,7) = "AGGACG"
name(8,6) = "A14844"
name(8,7) = "TGTGAA"
name(9,5) = "SISDA"
name(9,6) = "GCCACT"
name(0,7) = "TCR-CA1"
name(0,8) = "TGTGAA"
name(1,7) = "TCR-CB22"
name(1,8) = "GCTGCG"
name(2,7) = "TCR-CB22"
name(2,8) = "CTCTAA"
name(3,7) = "TCR-CB24"
name(3,8) = "GTCATG"
name(4,7) = "TCR-CB25"
name(4,8) = "GTTAGCC"
name(5,7) = "TCR-CB27"
name(5,8) = "ACCTTA"
name(6,7) = "VB12a"
name(6,8) = "ACGTCG"
name(7,7) = "VB12b"
name(7,8) = "CACTCA"
maxnum = 0
maxnum = 0

```


